
DATA FUNDAMENTALS

“Reports that say that something hasn't happened are always interesting to me, because as we know, there are **known knowns**; there are things we know that we know. There are **known unknowns**; that is to say, there are things that we now know we don't know. But there are also **unknown unknowns** – there are things we do not know we don't know.”

Donald Rumsfeld, US Department of Defense News Briefing, 2002

OUTLINE

1. Data 101 – Basic Data Concepts
2. Some Practical Definitions
3. Workflows and Pipelines – the Process of Working with Data
4. Models and Systems Thinking
5. Ethical Considerations and Best Practices

DATA 101 – BASIC DATA CONCEPTS

DATA FUNDAMENTALS

“You can have data without information, but you cannot have information without data.”

Daniel Keys Moran (attributed)

MODULE LEARNING OBJECTIVES

Preliminary familiarity with the following concepts:

- data, attribute (property, factor, variable)
- predictive models, explanatory models
- classification, class probability estimation, clustering, association rules, time series analysis, anomaly detection, decision tree, supervised learning, unsupervised learning

Compare and contrast: data science vs analytics (Business Intelligence).

Awareness of appropriate levels of trust in models.

WHAT IS DATA? WHERE DOES IT COME FROM?

4,529

'red'

25.782

'Y'

OBJECTS AND ATTRIBUTES



Object: apple

Shape: spherical

Colour: red

Function: food

Location: fridge

Owner: Jen

Remember: a person or an object is not simply the sum of its attributes!

FROM ATTRIBUTES TO DATASETS

Attributes are **fields** (or columns) in a database; objects are **instances** (or rows)

Objects are described by their **feature vector**, the collection of attributes associated with value(s) of interest

ID#	Shape	Colour	Function	Location	Owner
1	spherical	red	food	fridge	Jen
2	rectangle	brown	food	office	Pat
3	round	white	tell time	lounge	School
...

POISONOUS MUSHROOMS DATASET



Amanita muscaria

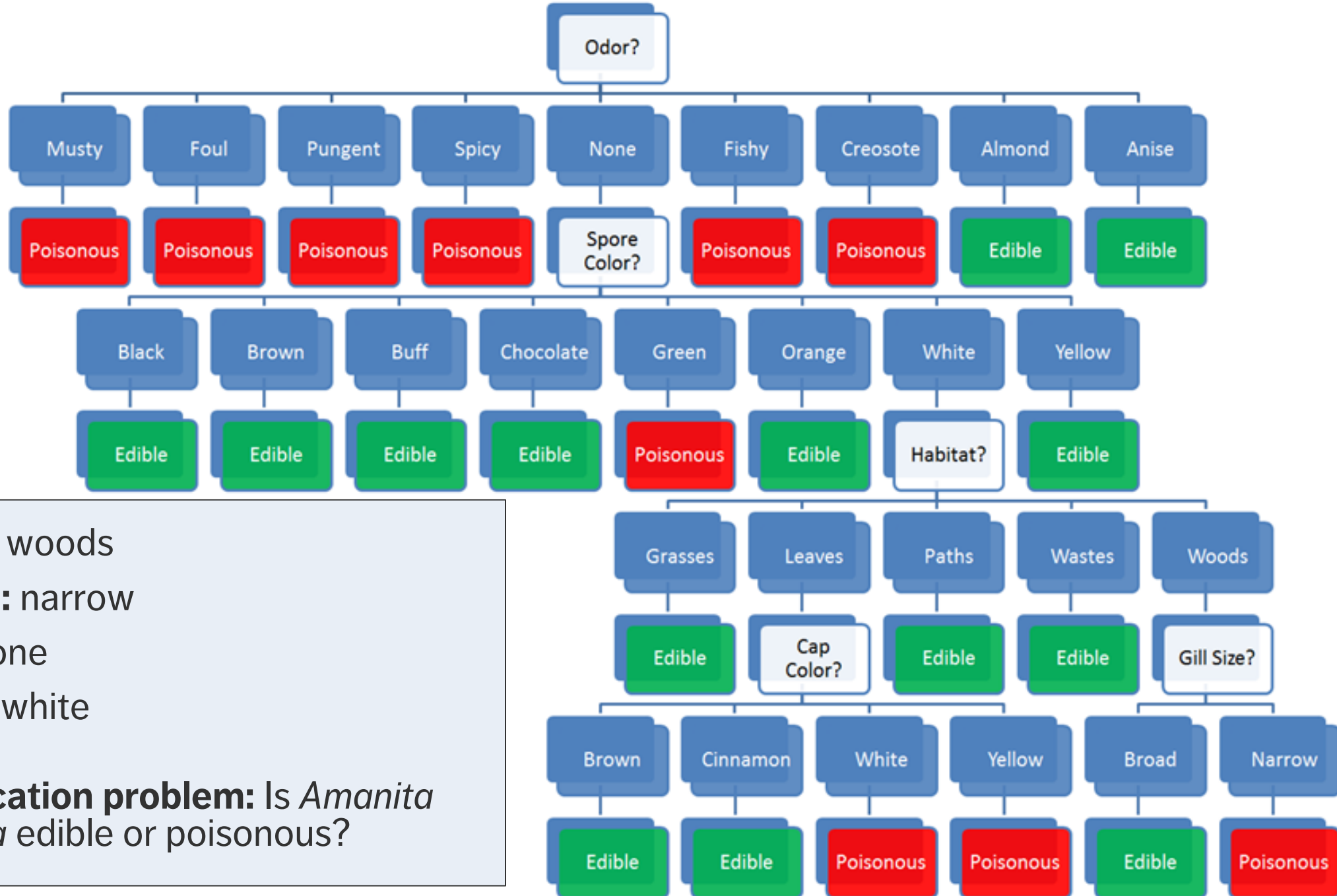
Habitat: woods

Gill Size: narrow

Odor: none

Spores: white

Classification problem: Is *Amanita muscaria* edible, or poisonous?



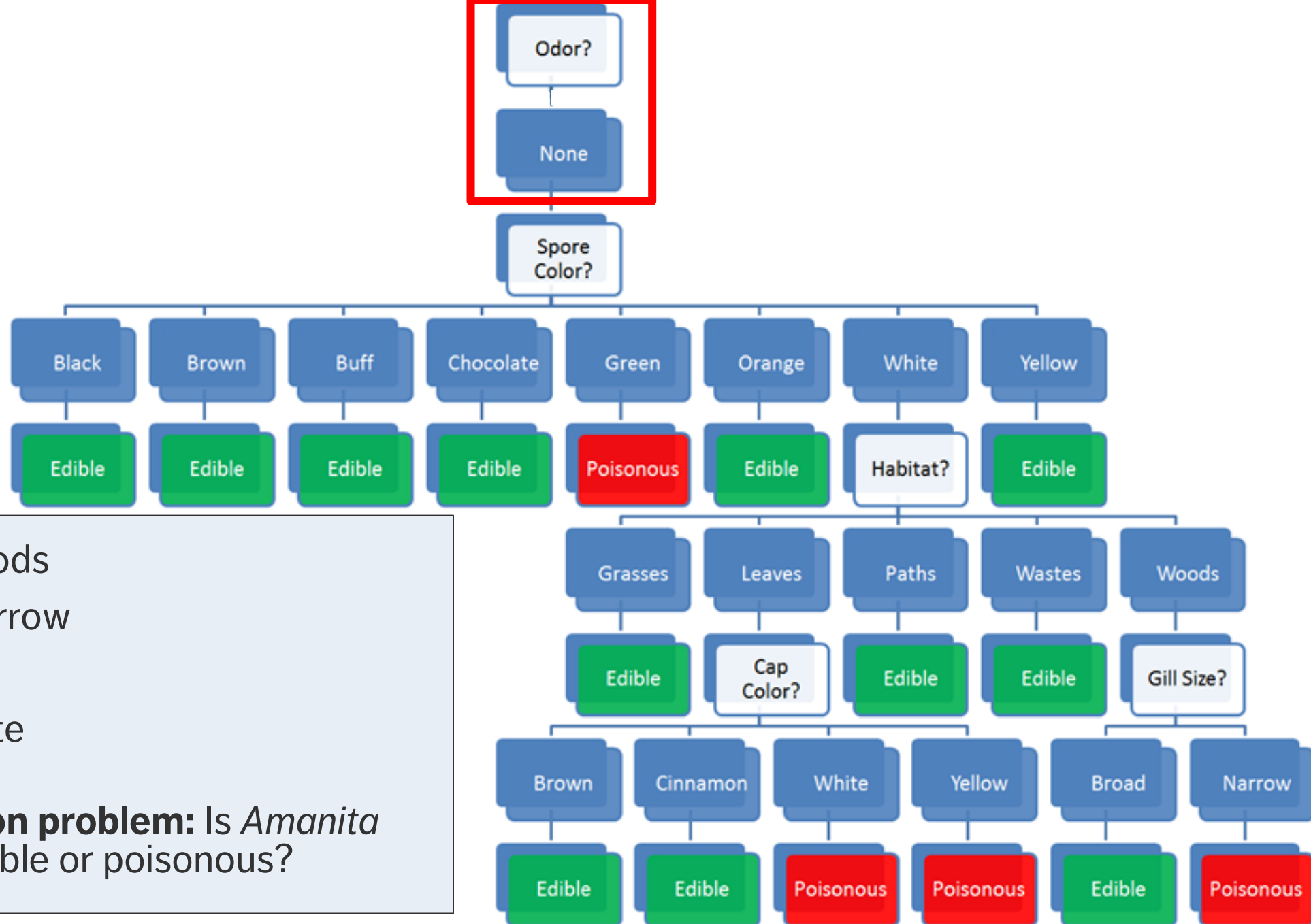
Habitat: woods

Gill Size: narrow

Odor: none

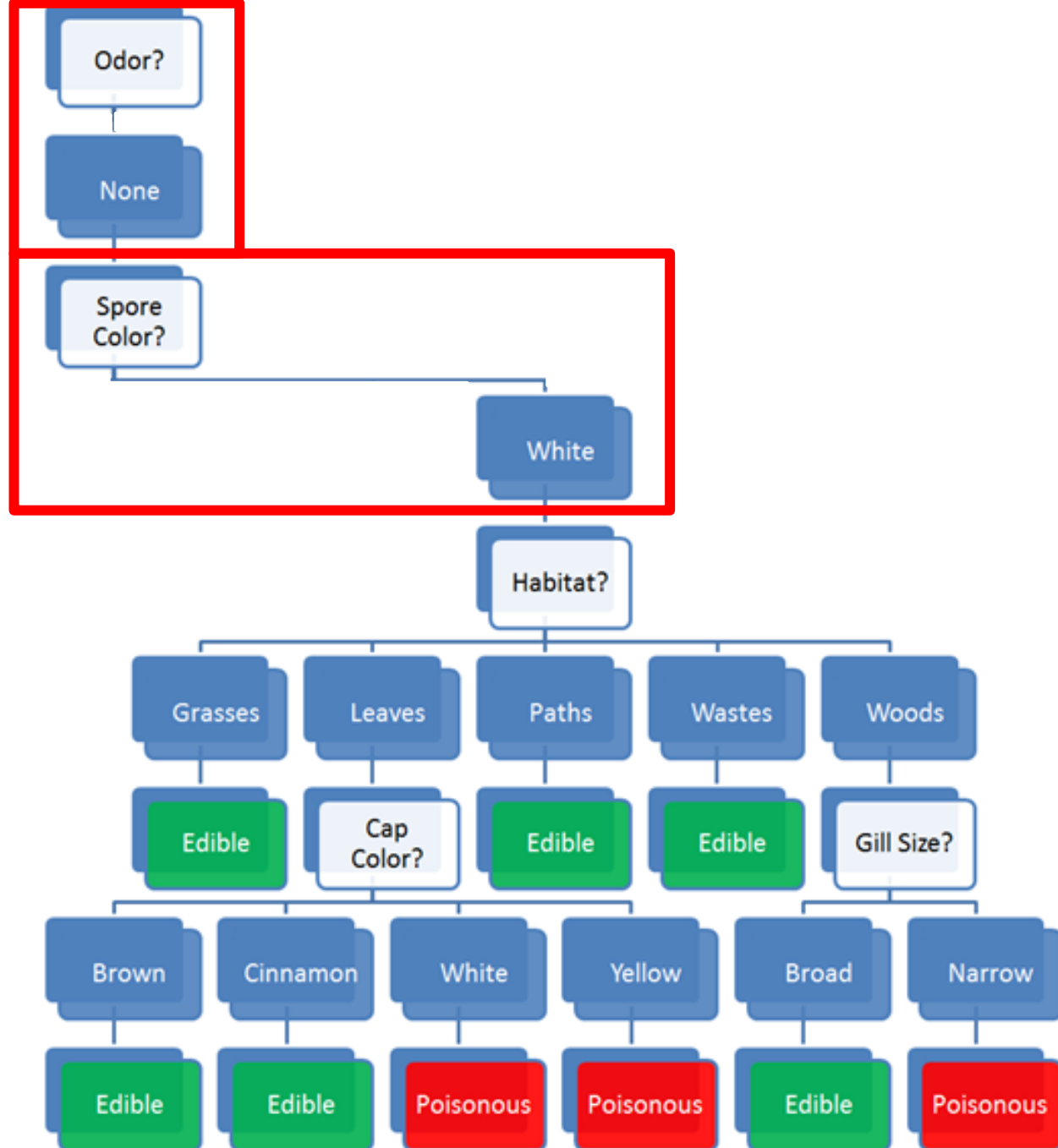
Spores: white

Classification problem: Is *Amanita muscaria* edible or poisonous?



Habitat: woods
Gill Size: narrow
Odor: none
Spores: white

Classification problem: Is *Amanita muscaria* edible or poisonous?



Habitat: woods

Gill Size: narrow

Odor: none

Spores: white

Classification problem: Is *Amanita muscaria* edible or poisonous?

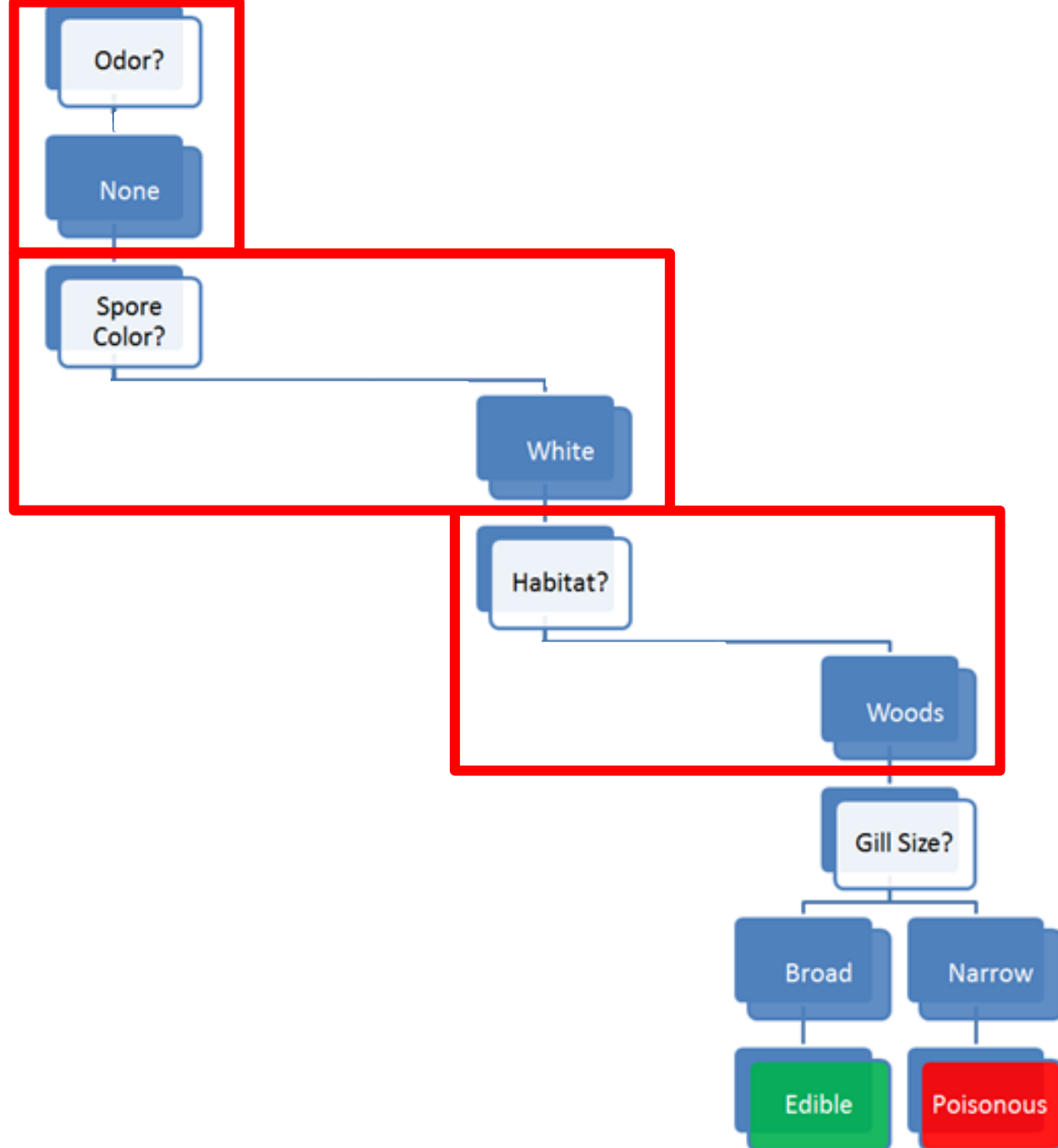
Habitat: woods

Gill Size: narrow

Odor: none

Spores: white

Classification problem: Is *Amanita muscaria* edible or poisonous?



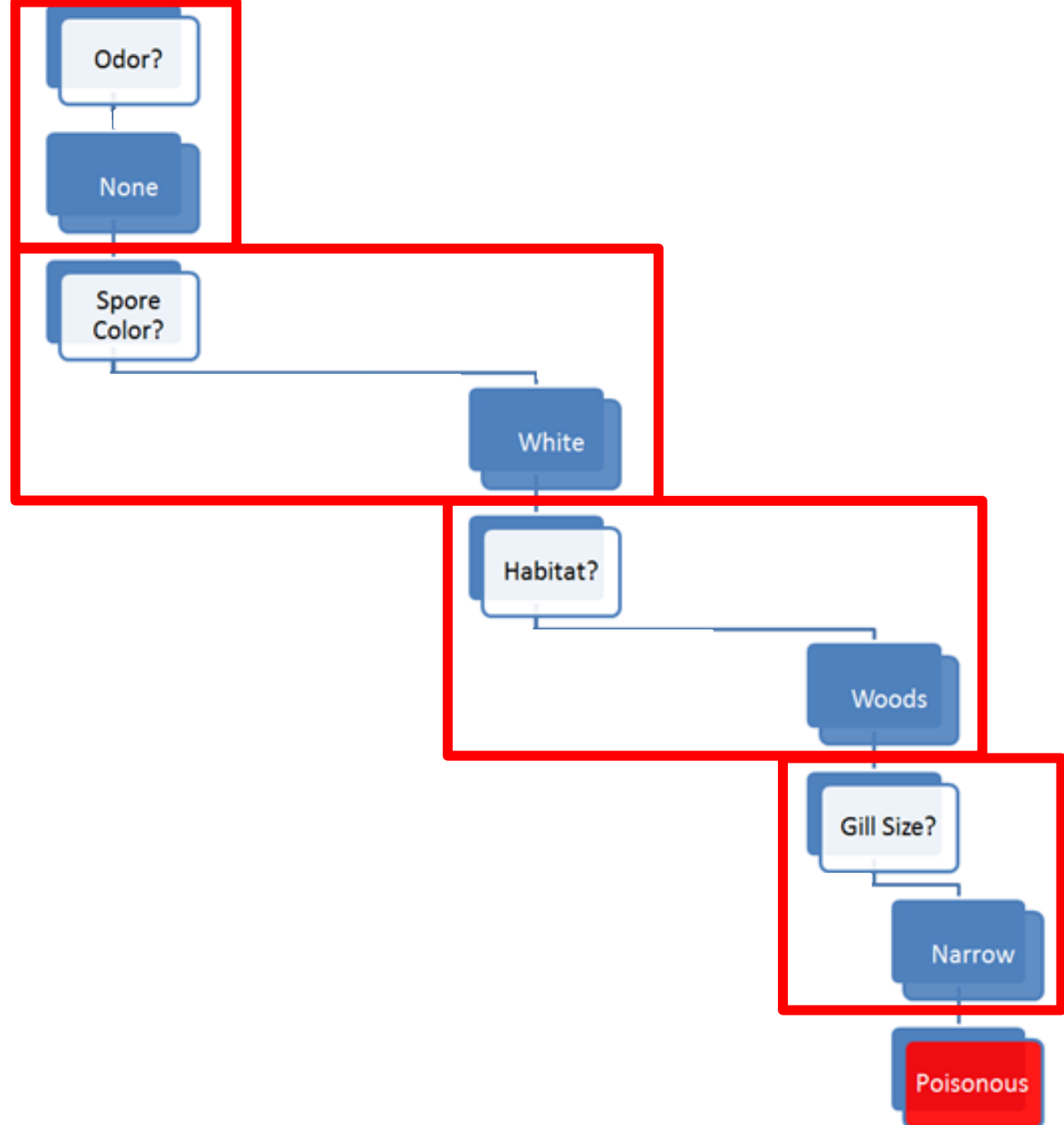
Habitat: woods

Gill Size: narrow

Odor: none

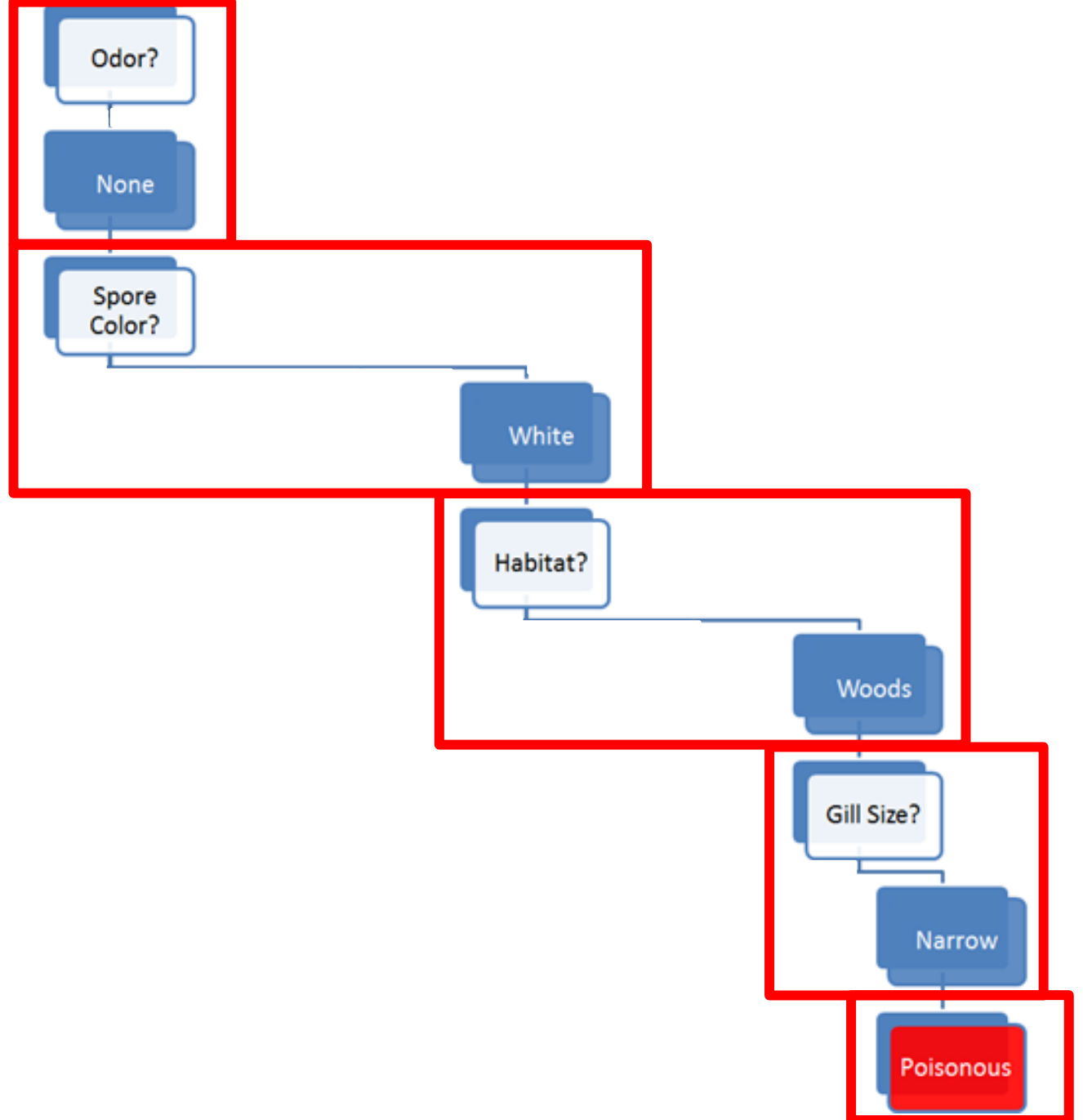
Spores: white

Classification problem: Is *Amanita muscaria* edible or poisonous?



Habitat: woods
Gill Size: narrow
Odor: none
Spores: white

Classification problem: Is *Amanita muscaria* edible or **poisonous**?



DISCUSSION

Would you trust an “**edible**” prediction?

Where is the model coming from?

What would you need to know to trust the model?

What’s the cost of making a classification mistake, in this case?

ASKING THE RIGHT QUESTIONS

Data science is really about asking and answering questions:

- **Analytics:** “How many clicks did this link get?”
- **Data Science:** “Based on this user’s previous purchasing history, can I predict what links they will click on the next time they access the site?”

Data mining/science models are usually **predictive** (not **explanatory**): they show connections, but don't reveal why these exist.

Warning: not every situation calls for data science, artificial intelligence, machine learning, or analytics.

DATA SCIENCE/MACHINE LEARNING/A.I. TASKS

Classification and class probability estimation: which clients are likely to be repeat customers?

Clustering: do customers form natural groups?

Association rule discovery: what books are commonly purchased together?

Others:

profiling and behaviour description; link prediction; value estimation (how much is a client likely to spend in a restaurant); **similarity matching** (which prospective clients are similar to a company's best clients?); **data reduction; influence/causal modeling**, etc.

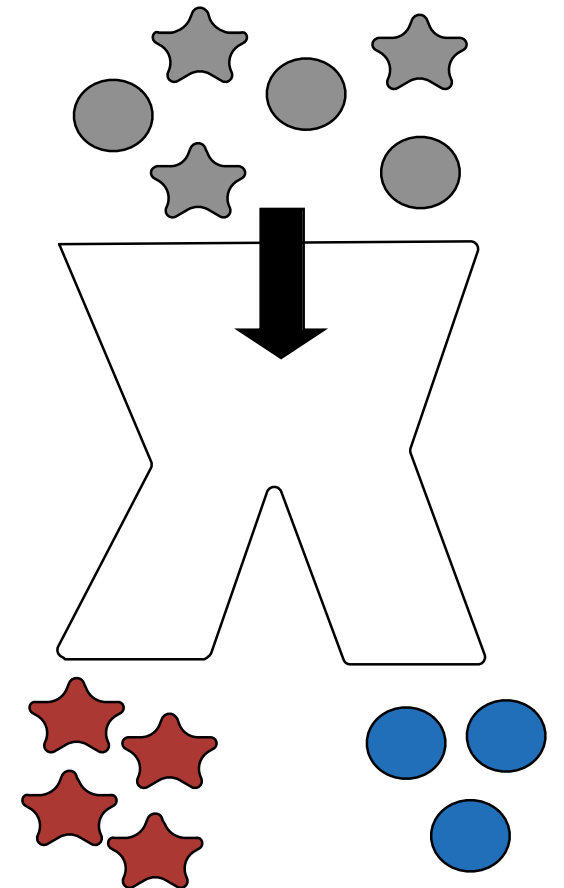
CLASSIFICATION

Classifier: If I'm presented with an object, can I classify it into one of several predefined categories?

Many different techniques to carry this out, but the steps are the same:

- Use a *training set* to teach the classifier how to classify.
- Test/validate the classifier using *new data*
- Use the classifier to classify *novel instances*

Some classifiers (e.g. neural nets) are very 'black box'. They might be good at classifying, but you don't know why!



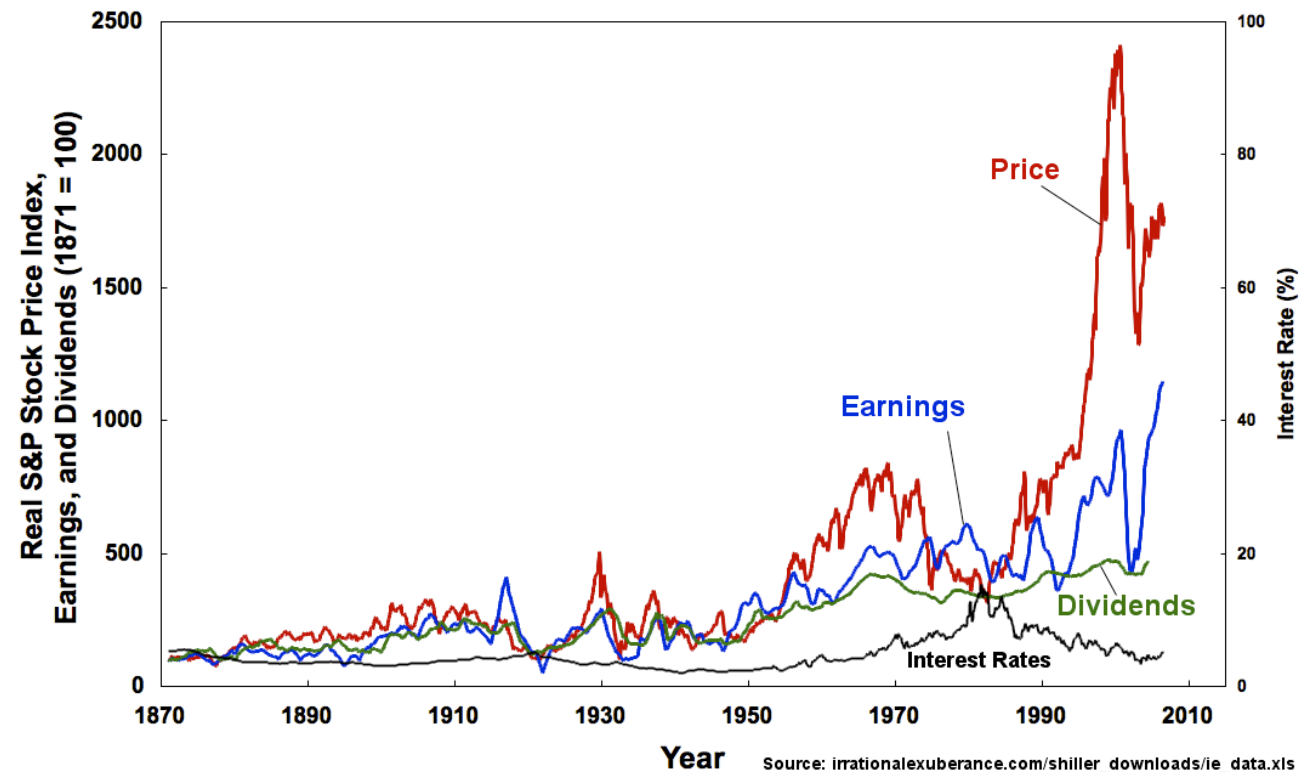
TIME SERIES ANALYSIS

A simple time series:

- Has two variables: time + 2nd variable
- The second variable is *sequential*

What is the pattern of behaviour of this second variable over time?
Relative to other variables?

Can we use this information to forecast the behaviour of the variable in the future?



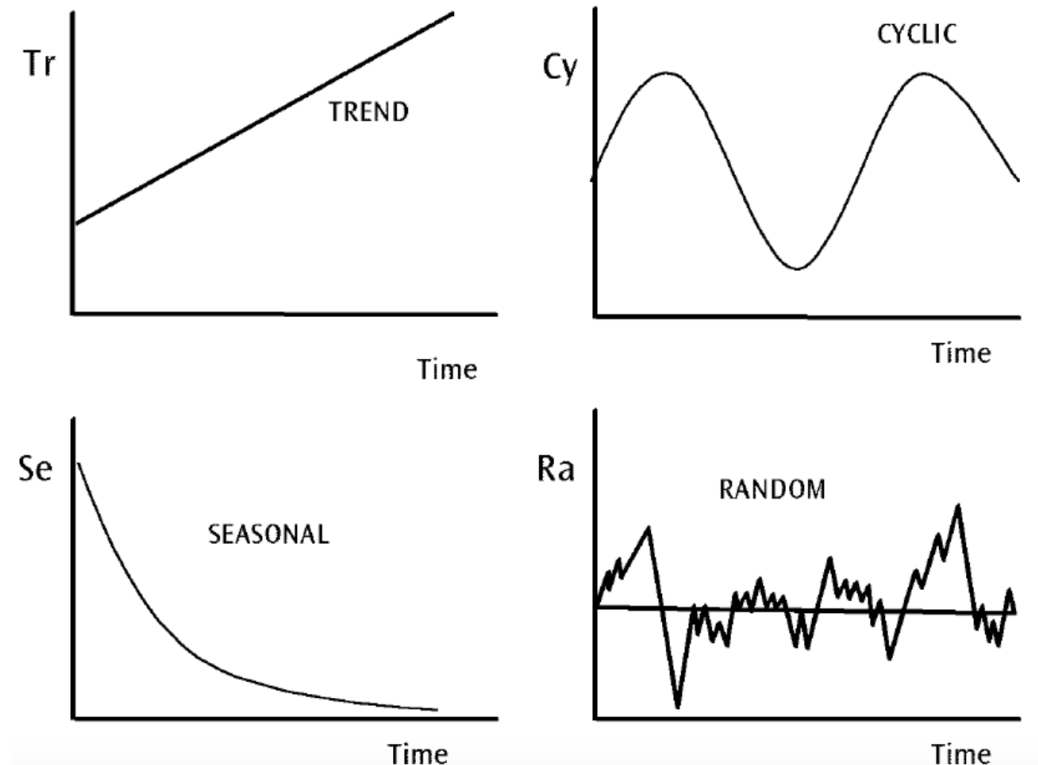
TEMPORAL PATTERNS

The goal here is our familiar analysis goals:

- find patterns in the data
- create a (mathematical) model that captures the essence of these patterns

The patterns can be quite complex – some fancy analysis typically required!

The data can often be broken down into multiple **component models**. There are software libraries that can help!



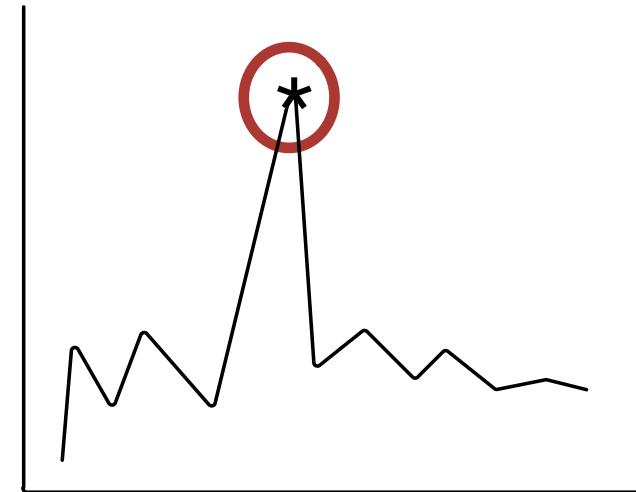
ANOMALY DETECTION

Anomaly: an unexpected, unusual, atypical or statistically unlikely event

Wouldn't it be nice to have a data analysis pipeline that alerted you when things were out of the ordinary?

Many different analytic approaches to take!

- Clustering
- Naïve Bayes
- Association rules deviation
- Ensemble techniques



ANOMALY DETECTION CASE STUDY

Energy 157 (2018) 336–352



ELSEVIER

Contents lists available at ScienceDirect

Energy

journal homepage: www.elsevier.com/locate/energy



Automated load pattern learning and anomaly detection for enhancing energy management in smart buildings

Alfonso Capozzoli*, Marco Savino Piscitelli, Silvio Brandi, Daniele Grassi, Gianfranco Chicco

Dipartimento Energia "Galileo Ferraris", Politecnico di Torino, Corso Duca degli Abruzzi 24, 10129, Torino, Italy

ARTICLE INFO

Article history:
Received 5 February 2018
Accepted 19 May 2018
Available online 21 May 2018

Keywords:
Energy consumption
Building energy management
Adaptive symbolic aggregate approximation
Anomaly detection
Data mining
Smart buildings

ABSTRACT

The energy management of buildings currently offers a powerful opportunity to enhance energy efficiency and reduce the mismatch between the actual and expected energy demand, which is often due to an anomalous operation of the equipment and control systems. In this context, the characterisation of energy consumption patterns over time is of fundamental importance. This paper proposes a novel methodology for the characterisation of energy time series in buildings and the identification of infrequent and unexpected energy patterns. The process is based on an enhanced Symbolic Aggregate approximation (SAX) process, and it includes an optimised tuning of the time window width and of the symbol intervals according to the building energy behaviour. The methodology has been tested on the whole electrical load of buildings for two case studies, and its flexibility and robustness have been confirmed. In order to demonstrate the implications for a preliminary diagnosis, some unexpected trends of the total electrical load have also been discussed in a post-mining phase, using additional datasets related to heating and cooling electrical energy needs.

The process can be used to support stakeholders in characterising building behaviour, to define appropriate energy management strategies, and to send timely alerts based on anomaly detection outcomes.

© 2018 Elsevier Ltd. All rights reserved.

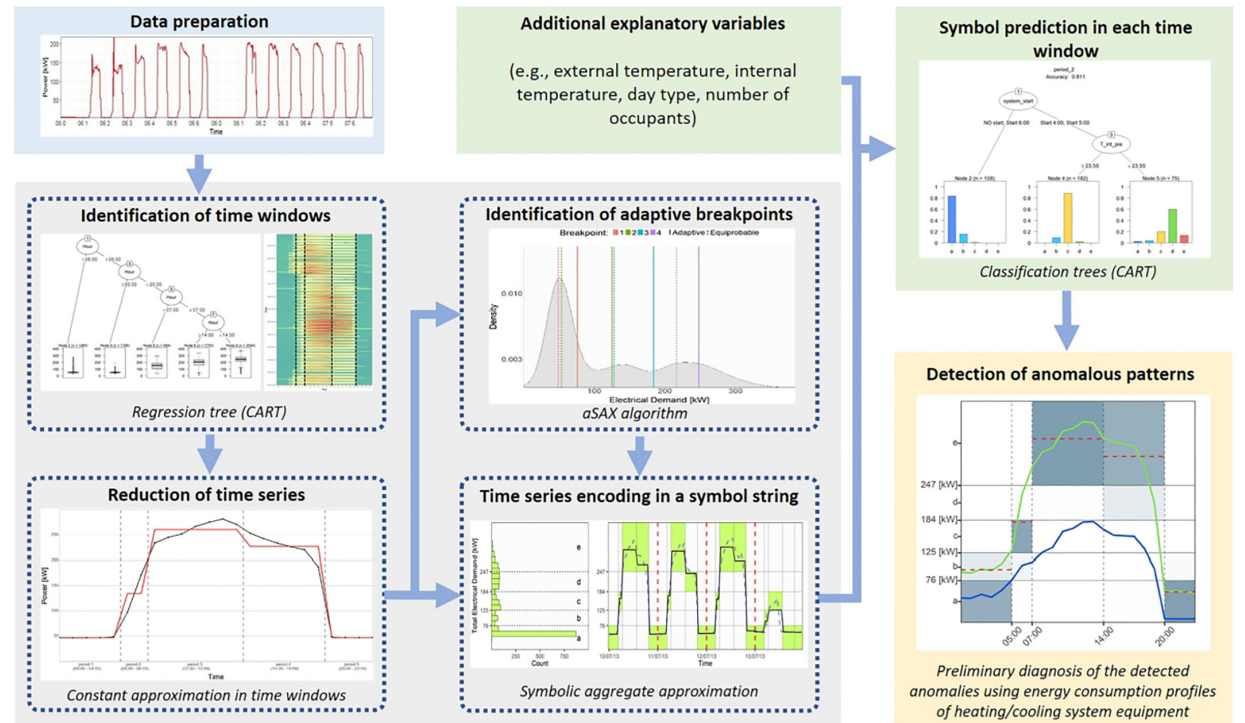


Fig. 2. – Framework for advanced energy consumption characterisation in buildings and anomalous pattern detection.

UNSUPERVISED LEARNING TECHNIQUES

Automated behaviours vs intelligent behaviours

Supervised: we give you some examples, you learn from them

Unsupervised: you learn on your own, based on what you experience

Unsupervised techniques:

- Association rules
- Recommender engines
- Novel categories (clustering)



SOME PRACTICAL DEFINITIONS

DATA FUNDAMENTALS

“What’s in a name? That which we call a rose
By any other name would smell as sweet.”

W. Shakespeare, Romeo and Juliet, Act II, Scene 2

MODULE LEARNING OBJECTIVES

Preliminary familiarity with the following concepts:

- data analysis
- data science
- machine learning
- patterns
- system
- artificial intelligence
- augmented intelligence

WHAT IS DATA ANALYSIS?

Finding **patterns** in data

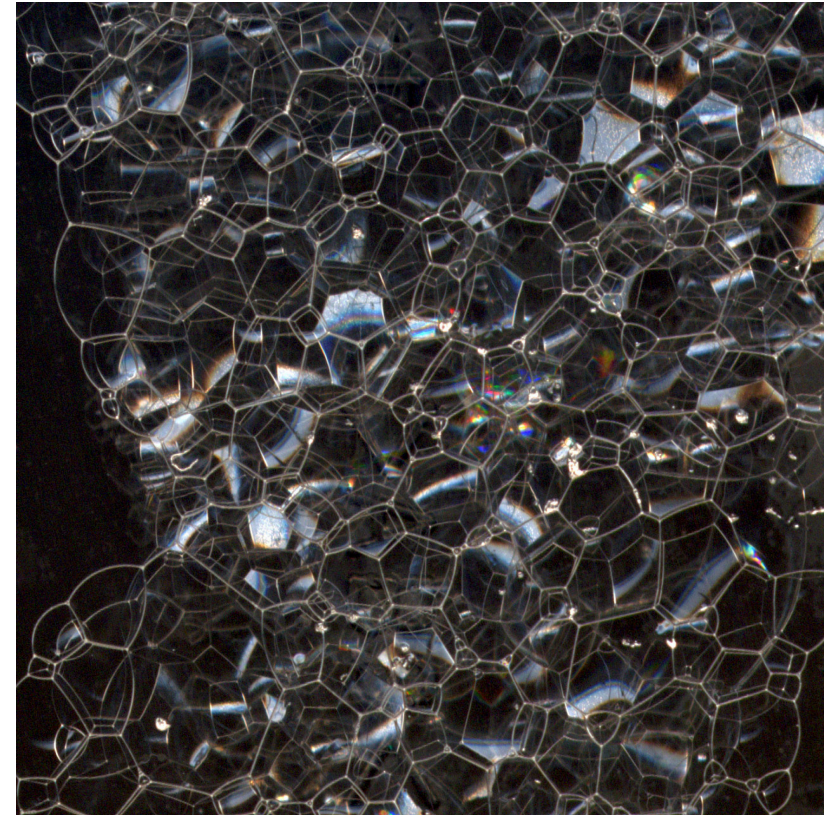
Using data to do something (answer a question, help decision-making, predict the future, draw a conclusion)

Creating models of your data

Describing or explaining your situation (your **system**)

(Testing (scientific) hypotheses?)

(Carrying out calculations on data?)



The more complicated the pattern, the more complicated the analysis (?)

WHAT IS DATA SCIENCE?

Data science is the collection of processes by which we extract useful and **actionable insights** from data.

T. Kwartler (paraphrased)

Data science is the **working intersection** of statistics, engineering, computer science, domain expertise, and “hacking.” It involves two main thrusts: **analytics** (counting things) and **inventing new techniques** to draw insights from data.

H. Mason (paraphrased)

WHAT IS MACHINE LEARNING?

Starting around the 1940s researchers began in earnest to teach machines how to learn

The goal of **machine learning** was to create machines that could learn and adapt and respond to novel situations

A wide variety of techniques, accompanied by a great deal of theoretical underpinning, was created in an effort to achieve this goal



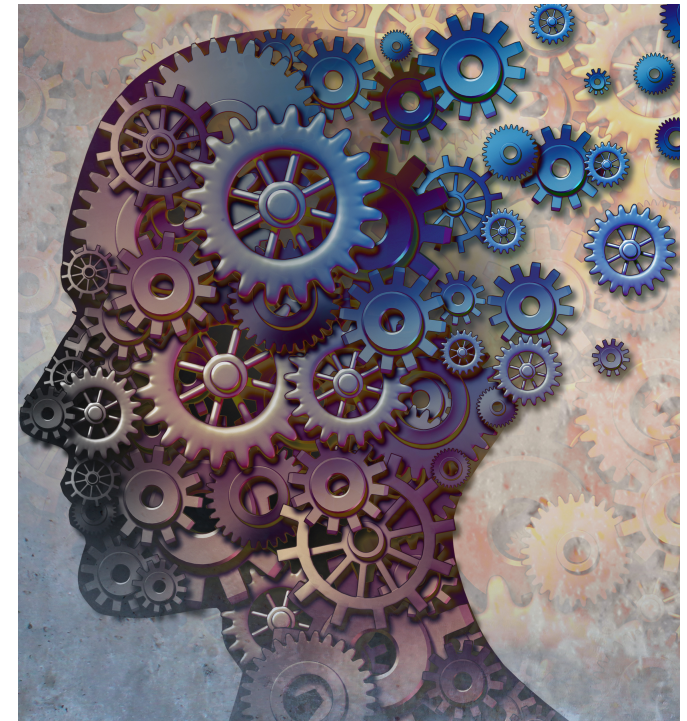
WHAT IS ARTIFICIAL/AUGMENTED INTELLIGENCE?

Artificial Intelligence (A.I.) is non-human intelligence that has been engineered rather than one that has evolved naturally.

Artificial intelligence research is research carried out in pursuit of this goal.

Pragmatically speaking, A.I. is “computers carrying out tasks that only humans can usually do”.

Augmented Intelligence is human intelligence that is supported or enhanced by machine intelligence.



WORKFLOWS AND PIPELINES

DATA FUNDAMENTALS

“All models are wrong. Some models are useful.”

George Box

MODULE LEARNING OBJECTIVES

Preliminary familiarity with the following concepts:

- workflow and components (data collection, data exploration, etc.)
- analytical model
- data mining
- analytic decay
- data science ecosystem
- data science teams

Awareness of the non-linearity of the data analytical process.

THE DATA SCIENCE “WORKFLOW”

Objective/
Rationale

Data
Collection

Data
Exploration

Utilization and
Decision
Support

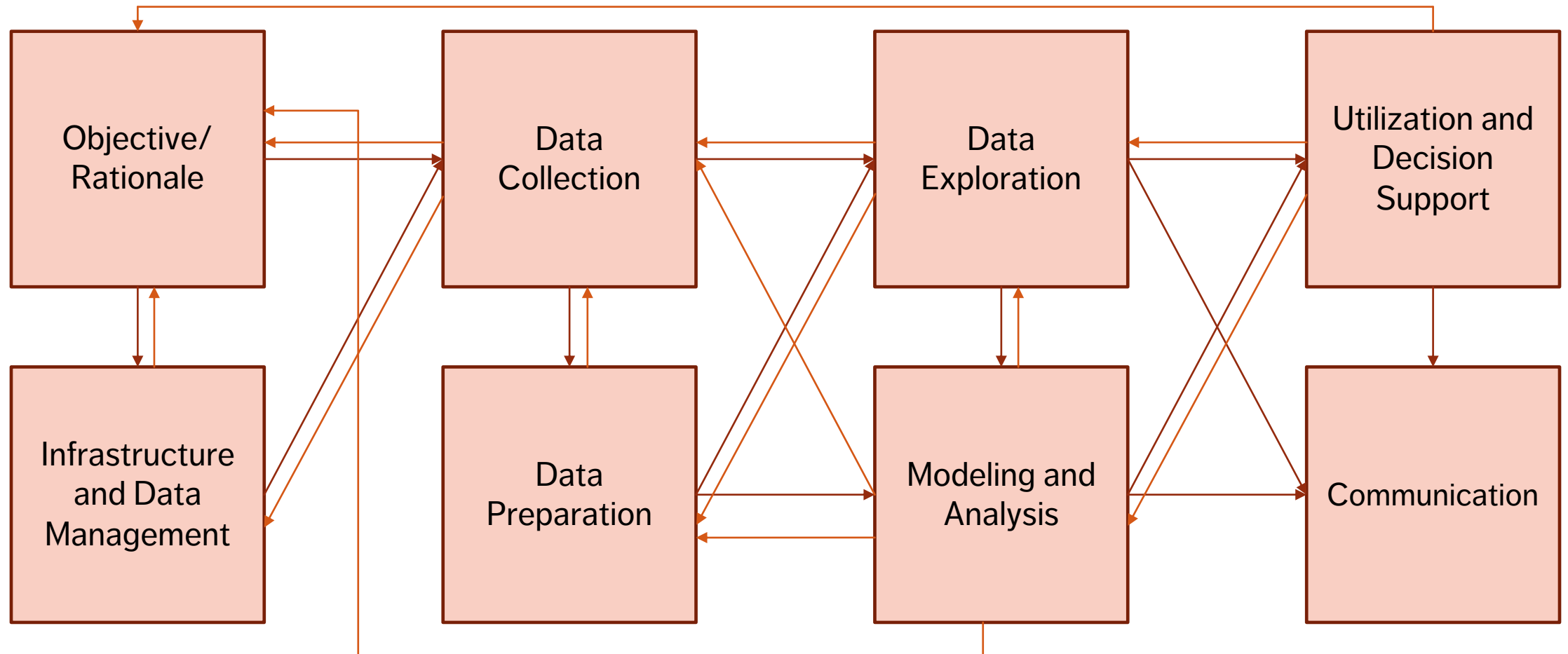
Infrastructure
and Data
Management

Data
Preparation

Modeling and
Analysis

Communication

THE DATA SCIENCE “WORKFLOW”



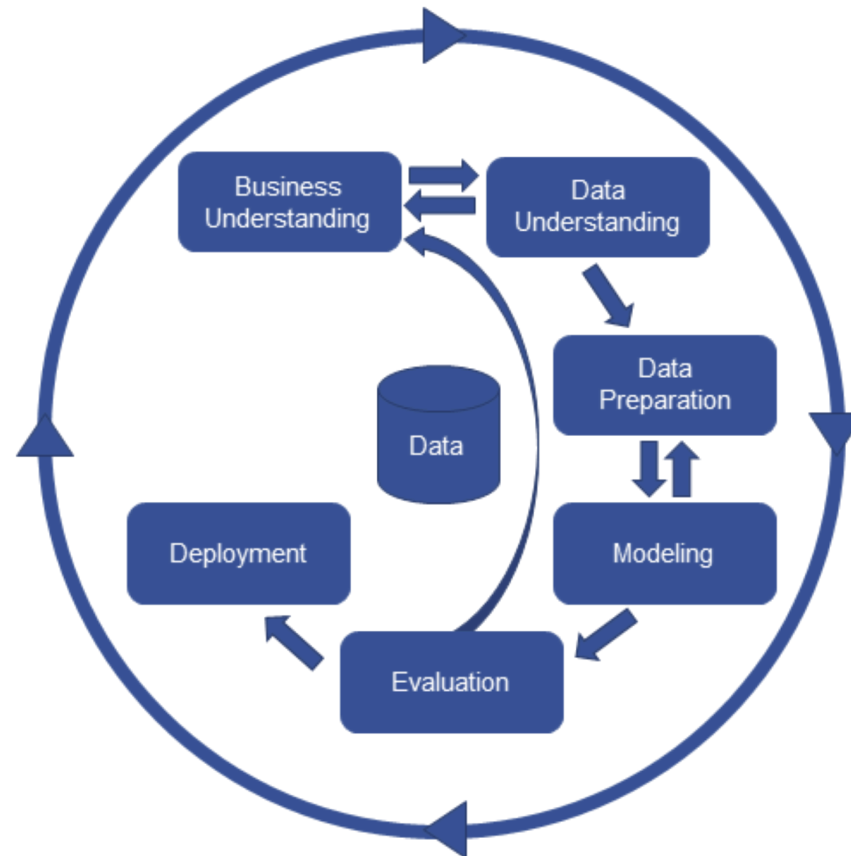
THE DATA ANALYSIS PROCESS

A **large number of analytical models** have to be generated before a final selection can be made.

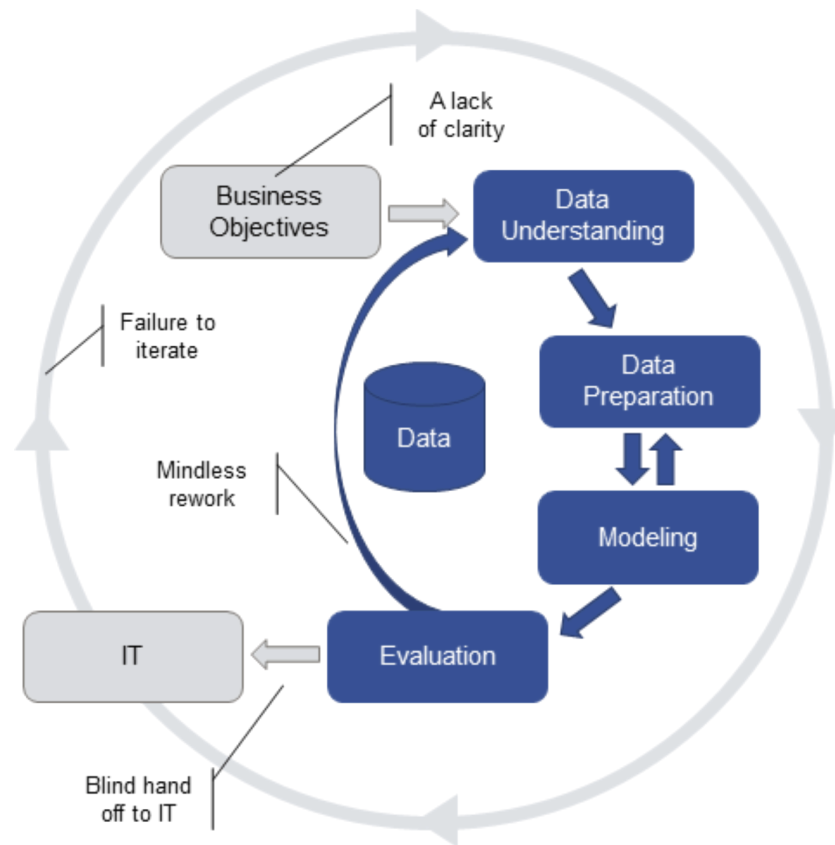
Iterative process: feature selection and data reduction may require numerous visits to domain experts before models start yielding promising results.

Domain-specific knowledge has to be integrated in the models in order to beat random classifiers and clustering schemes, **on average**.

CROSS INDUSTRY STANDARD PROCESS, DATA MINING



CROSS INDUSTRY STANDARD PROCESS, DATA MINING



LIFE AFTER ANALYSIS

When an analysis or model is ‘released into the wild’, it can take on a life of its own.

Analysts may eventually have to relinquish control over dissemination. Results may be misappropriated, misunderstood, or shelved. What can be done to prevent this?

Because of **analytic decay**, better to see the last analytical step NOT as a static dead end, but rather as an **invitation to return to the beginning of the process**.

DATA SCIENCE ECOSYSTEM

Data analysis is a **team sport**, with team members needing a good understanding of both **data** and **context**

- data management
- data preparation
- analysis
- communications

Even slight improvements over a current approach can find a useful place in an organization – **data science is not solely about Big Data and disruption!**

MODELS AND SYSTEMS THINKING

DATA FUNDAMENTALS

“What if the only valid model of the Universe is the Universe itself?”

Unknown

MODULE LEARNING OBJECTIVES

Preliminary familiarity with the following concepts:

- representation
- systems
- models
- properties
- knowledge gap
- conceptual model

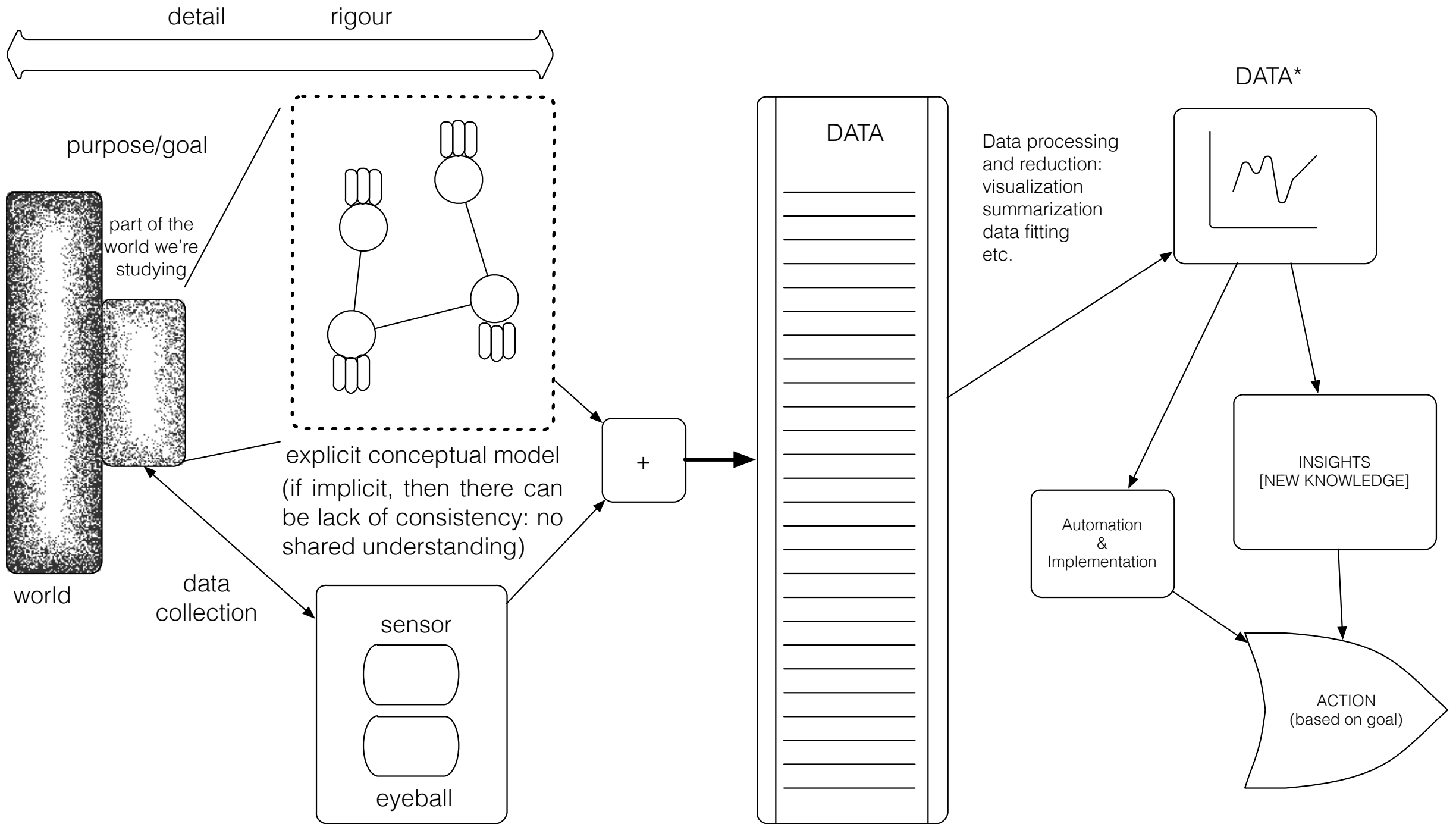
REPRESENTATION

A representation is an object that stands in for another object.

A representation may or may not physically resemble the object it represents.

Representations of the world help us to understand, navigate and manipulate the world.





THINKING IN SYSTEMS TERMS

In order to understand how various aspects of the World interact with one another, we need to **carve out chunks** corresponding to the aspects and define their **boundaries**.

Working with other intelligences requires **shared understanding** of what is being studied.

A **system** is made up of **objects** with **properties** that potentially change over time. Within the system we perceive **actions** and **evolving** properties leading us to think in terms of **processes**.

THINKING IN SYSTEMS TERMS

Objects themselves have various properties. Natural processes generate (or destroy) objects, and may change the properties of these objects over time.

We **observe**, **quantify**, and **record** particular values of these properties at particular points in time.

This generates data points, capturing the **underlying reality** to some degree of **accuracy** and **error** (biased or unbiased).

IDENTIFYING GAPS IN KNOWLEDGE

A **gap in knowledge** is identified when we realize that what we thought we knew about a system proves incomplete (or false).

This might happen repeatedly, at any moment in the process:

- data cleaning
- data consolidation
- data analysis

The solution is to be flexible. When faced with such a gap, **go back, ask questions,** and **modify the system representation.**

CONCEPTUAL MODELS

Exercise:

- assume that an acquaintance has just set foot in your living space for the first time.
- you are on the phone with them but not currently at home.
- explain to them how to go about preparing a cup of sugar.

Conceptual models are built using methodical investigation tools

- diagrams
- structured interviews
- structured descriptions
- etc.

RELATING THE DATA TO THE SYSTEM

Is the data which has been collected and analyzed going to be of any use when it comes to understanding the system?

This question can only be answered if we understand:

- how the data is **collected**
- the **approximate nature** of both data and system
- what the data **represents** (observations and features)

Is the combination of system and data **sufficient** to understand the aspects of the world under consideration?

TAKE-AWAYS

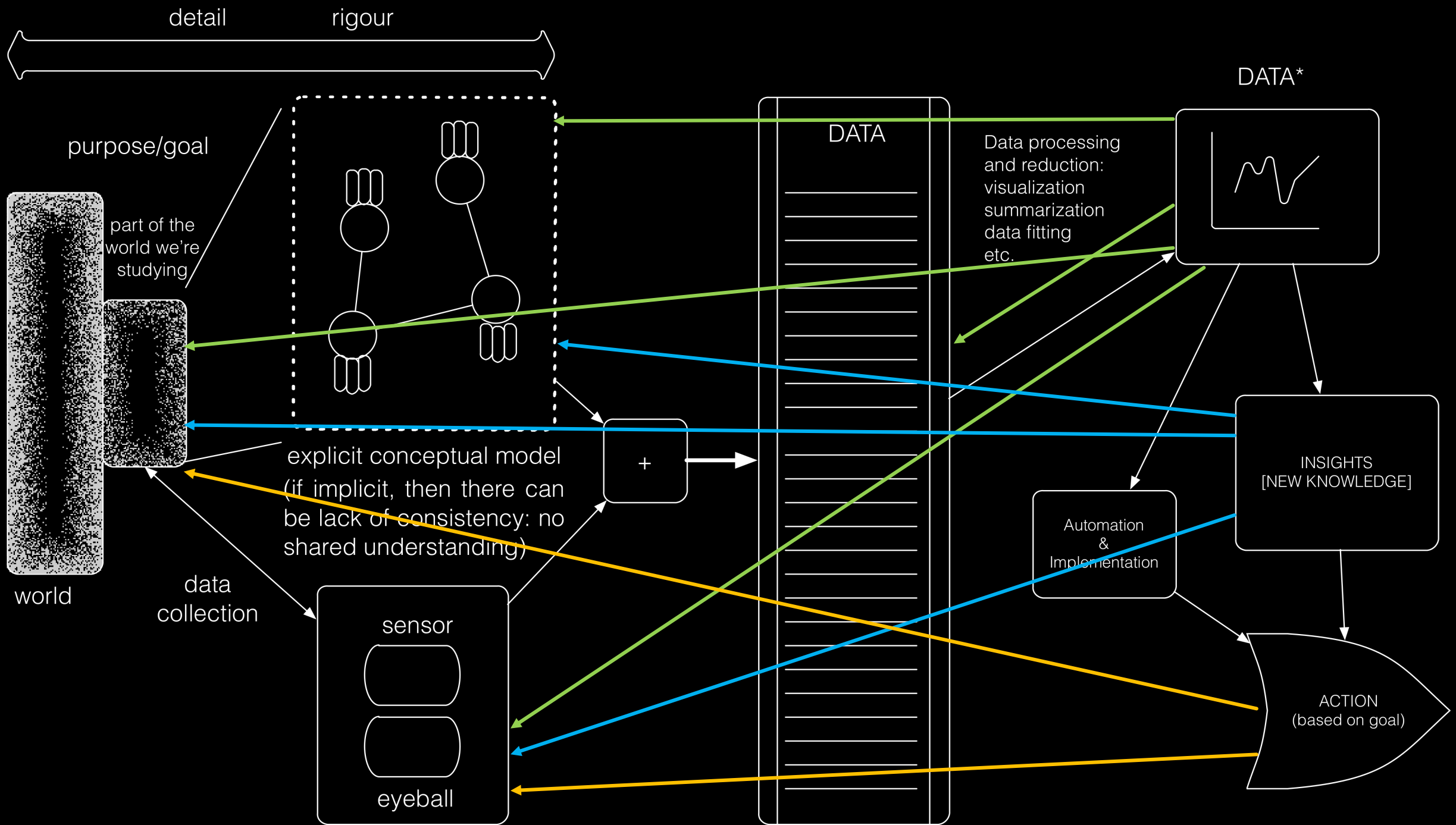
Certain aspects of the Universe can be approximated with the help of systems.

System models provide the basis under which data is identified and collected, but data itself is approximate and selective.

Knowledge gaps happen. Be prepared and ready to re-visit your set-up regularly.

We often only rely on implicit conceptual modeling, but there's danger that way.

If the data, the system, and the world are out of alignment, insights might prove useless.



ETHICAL CONSIDERATIONS AND BEST PRACTICES

DATA FUNDAMENTALS

“We have flown the air like birds and swum the sea like fishes, but have yet to learn the simple act of walking the Earth like brothers.”

Martin Luther King, Jr.

MODULE LEARNING OBJECTIVES

Preliminary familiarity with the following concepts:

- ethics and best practices
- First Nations principles (OCAP)
- “do no harm“
- informed consent
- privacy
- model validity

DISCUSSION

What harm can come from data?

THE NEED FOR ETHICS

Formerly: “**Wild West**” mentality to data collection (and use). Whatever wasn’t technologically forbidden was allowed.

Now: professional codes of conduct are being devised for data scientists (outline responsible ways to practice data science).

Additional responsibility for data scientists; but also **protection** against being hired to carry out questionable analyses.

Does your organization have a code of ethics for its data scientists? For its employees?

WHAT ARE ETHICS?

Broadly speaking, ethics refers to the **study** and **definition** of **right and wrong conducts**:

- “not [...] social convention, religious beliefs, or laws”. (R.W. Paul, L. Elder)

Influential *Western* ethical theories:

- Kant's **golden rule** (do unto others...), **consequentialism** (the ends justify the means), **utilitarianism** (act in order to maximize positive effect), etc.

Influential *Eastern* ethical theories:

- **Confucianism, Taoism, Buddhism** (?), etc.

WHAT ARE ETHICS?

First Nations Principles of **OCAP**®:

- **Ownership**
cultural knowledge, data, and information is owned by First Nations communities
- **Control**
First Nations communities have the right to control all aspects of research and information management that impact them
- **Access**
First Nations communities must have access to information and data about themselves no matter where it is held
- **Possession**
First Nations communities must have physical control of relevant data

ETHICS IN THE DATA CONTEXT

Data ethics questions:

- **Who**, if anyone, owns data?
- Are there **limits** to how data can be used?
- Are there **value-biases** built into certain analytics?
- Are there categories that should **not** be used in analyzing personal data?
- Should some data be **publicly available** to **all** researchers?

Analytically, the **general** is preferred to the **anecdotal** – decisions made on the basis of machine learning and A.I. (security, financial, marketing, etc.) may affect real beings in **unpredictable ways**.

BEST PRACTICES

“Do No Harm”: data collected from an individual **should not be used to harm** the individual.

Informed Consent:

- Individuals must **agree to the collection and use** of their data
- Individuals must have a **real understanding of what they are consenting to**, and of **possible consequences** for them and others

Respect “Privacy”: excessively hard to maintain in the age of constant trawling of the Internet for personal data.

BEST PRACTICES

Keep Data Public: data should be kept **public** (all? most? any?).

Opt-In/Opt-Out: Informed consent requires the ability to **opt out**.

Anonymize Data: removal of id fields from data prior to analysis.

“Let the Data Speak”:

- no cherry picking
- importance of validation (more on this later)
- correlation and causation (more on this later, too)
- repeatability

MODEL ASSESSMENT AND VALIDITY

Models should be **current**, **useful**, and **valid**.

Data can be used in conjunction with existing models to come to some conclusions, or can be used to update the model itself.

At what point does one determine that the current data model is **out-of-date** or is **not useful anymore**?

Past successes can lead to **reluctance** to re-assess and re-evaluate a model.

READINGS AND REFERENCES

DATA FUNDAMENTALS

REFERENCES

[First Nations – OCAP](#)

Wikipedia article on [Semi-Supervised Learning](#)

Wikipedia article on [Supervised Learning](#)

Wikipedia article on [Reinforcement Learning](#)

Wikipedia article on [Unsupervised Learning](#)

J. Blitzen [2017], [What is it like to design a data science class?](#), answer on Quora

J. Taylor [2017], [4 Problems with CRISP-DM](#), KDNuggets.

Brin, D. [1998], [The Transparent Society: Will Technology Force Us to Choose Between Privacy and Freedom?](#), Perseus.

REFERENCES

Mayer-Schönberger, V. and Cukier, K. [2013], [Big Data: A Revolution That Will Transform How We Live, Work, and Think](#), Eamon Dolan/Houghton Mifflin Harcourt.

Mayer-Schönberger, V. [2009], [Delete: The Virtue of Forgetting in the Digital Age](#), Princeton University Press.

Data Science Association, [Data Science Code of Professional Conduct](#).

Chen, M. [2013], [Is 'Big Data' Actually Reinforcing Social Inequalities?](#), The Nation.

Shin, L. [2013], [How the New Field of Data Science is Grappling With Ethics](#), SmartPlanet.

Schutt, R. and O'Neill, C. [2013], [Doing Data Science: Straight Talk From the Front Line](#), O'Reilly.

O'Neill, C. [2016], [Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy](#), Crown.

REFERENCES

Chang, R.M., Kauffman, R.J., Kwon, Y. [2014], *Understanding the paradigm shift to computational social science in the presence of big data*, Decision Support Systems, 63:67–80, Elsevier.

Hurlburt, G.F., Voas, J. [2014], *Big Data, Networked Worlds*, IEEE Computer Society.

Introna, L.D. [2007], *Maintaining the reversibility of foldings: Making the ethics (politics) of information technology visible*, Ethics and Information Technology, 9:11–25, Springer.

Floridi, L. [2011], *The philosophy of information*, Oxford University Press.

Floridi, L. (ed) [2006], *The Cambridge handbook of information and computer ethics*, Cambridge University Press, 2006.

[Big Data & Ethics](#)

Mason, H. [2012], [What is a Data Scientist?](#), Forbes.

REFERENCES

Schlimmer, J.S. [1987], *Concept Acquisition Through Representational Adjustment (Technical Report 87-19)*. Department of Information and Computer Science, UCalifornia, Irvine.

Iba, W., Wogulis, J., Langley, P. [1988], *Trading off Simplicity and Coverage in Incremental Concept Learning*, in Proceedings of the 5th International Conference on Machine Learning, 73-79. Ann Arbor, Michigan: Morgan Kaufmann.

Gorelik, B. [2017], [Don't study data science as a career move; you'll waste your time!](http://gorelik.net), gorelik.net

J. Leskovec, A. Rajaraman, J. Ullman [2015] *Mining of Massive Datasets*, Cambridge University Press.

Hastie, T., Tibshirani, R., and J. Friedman [2008], *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed., Springer.

Provost, F., Fawcett, T. [2013], *Data Science for Business*, O'Reilly.